# From Text to Graph

*Automatic knowledge extraction and semantification of texts*

60th CIDOC CRM SIG Meeting
03.04.2025

Stephen Hart, Unibe

# Knowledge Extraction: The Goal of the KNEX tool
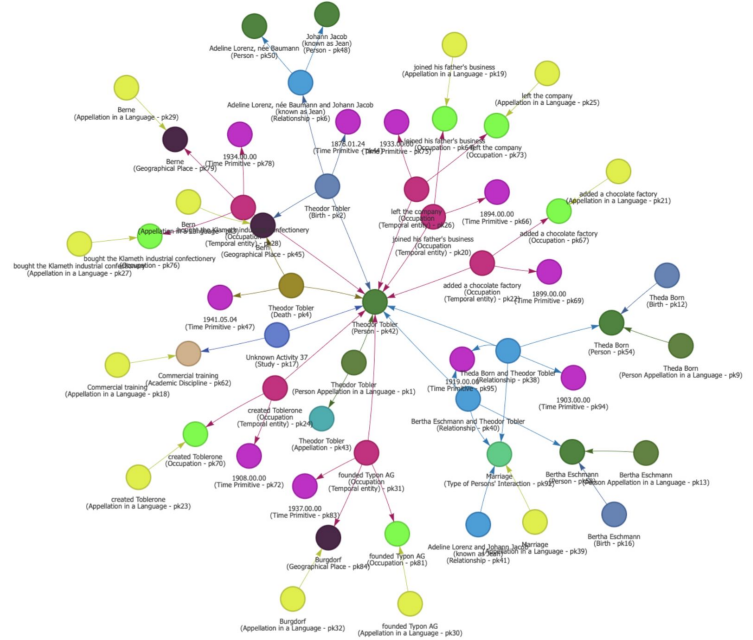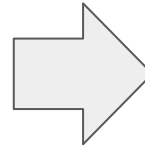
DE **FR** IT

## Theodor Tobler

## Daniel Peter

DE **FR** IT

## François-Louis Cailler

Version du: 20.03.2003

Auteure/Auteur: Gilbert Marion

∗ 11.6.1796 à Vevey, † 6.4.1852 à Corsier-sur-Vevey, prot., de Daillens et Vevey. Fils de François Louis. Beau-père de Daniel Peter. ∞ Louise Albertine Perret, de Boudry. Après un apprentissage chez un épicier de Vevey, C. séjourne en Italie du Nord où il découvre les chocolatiers tessinois de Turin et s'initie à la fabrication du chocolat. De retour à Vevey en 1818, il crée des machines pour ce qui devient en 1819 la première fabrique moderne de chocolat en Suisse, installée au lieudit En Copet, sur la commune de Corsier. Il fut l'initiateur de la présentation du chocolat sous forme de plaque.
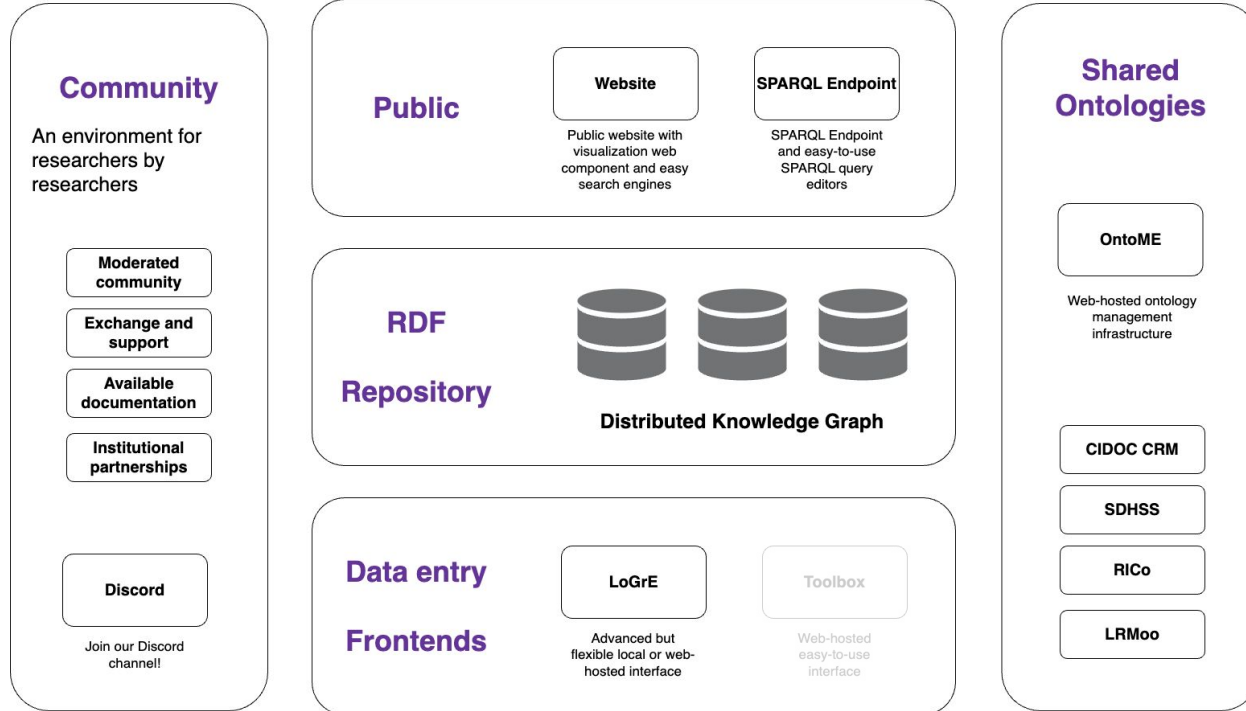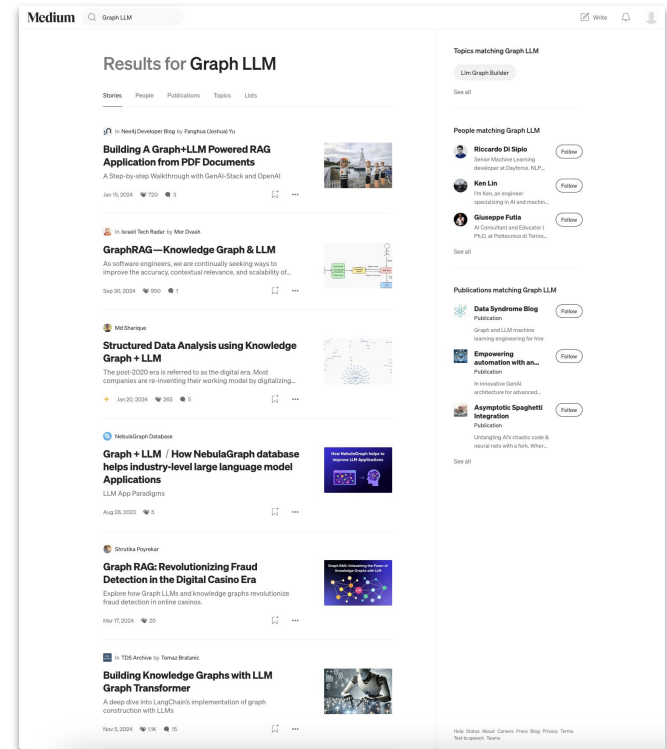
# Knex Context: The Geovistory Environment



**https://www.geovistory.org**

# The Geovistory Environment



**Community**

An environment for researchers by researchers

Moderated community

Exchange and support

Available documentation

Institutional partnerships

Discord

Join our Discord channel!

**Public**

Website

Public website with visualization web component and easy search engines

SPARQL Endpoint

SPARQL Endpoint and easy-to-use SPARQL query editors

**RDF Repository**

Distributed Knowledge Graph

**Data entry Frontends**

LoGrE

Advanced but flexible local or web-hosted interface

Toolbox

Web-hosted easy-to-use interface

**Shared Ontologies**

OntoME

Web-hosted ontology management infrastructure

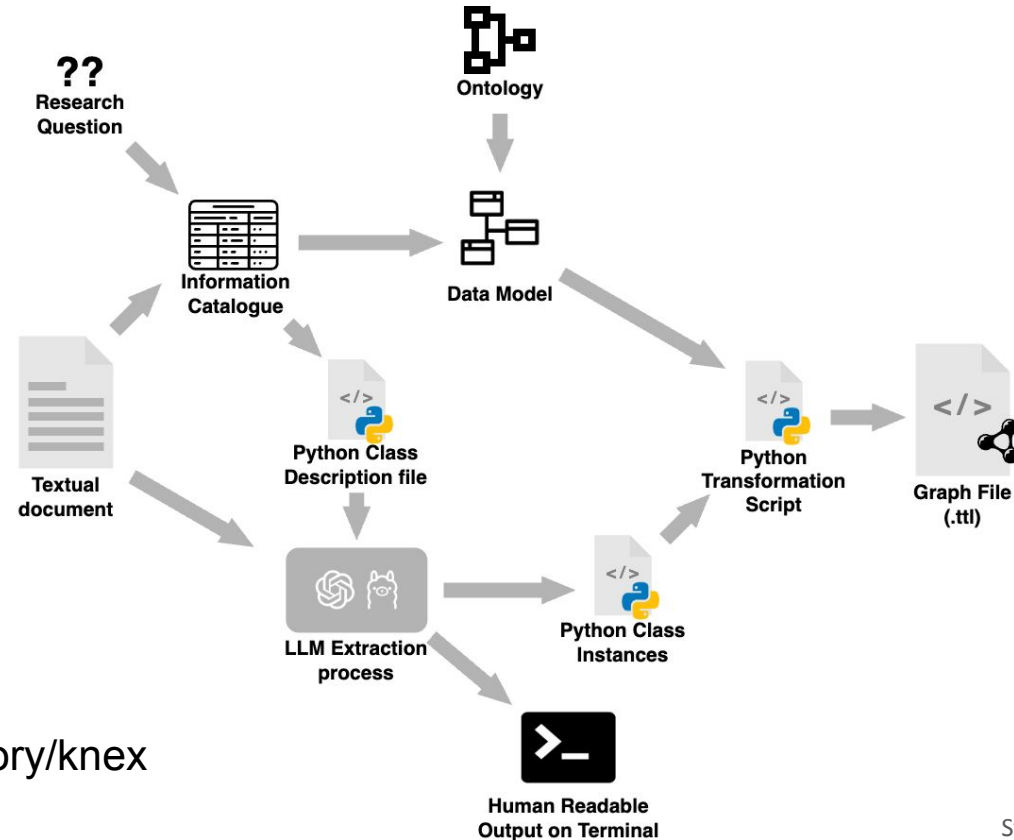CIDOC CRM

SDHSS

RICo

LRMoo

# Knowledge Extraction: The Goal

- New LLM technologies has transformed the way of interacting with texts and generating new knowledge
  - Many new attempts at generating Knowledge Graphs from texts
- For the moment little attempt in the Humanities (but it's starting)
- Solutions existing lack Semantic Engineering, and creates the ontology from the text
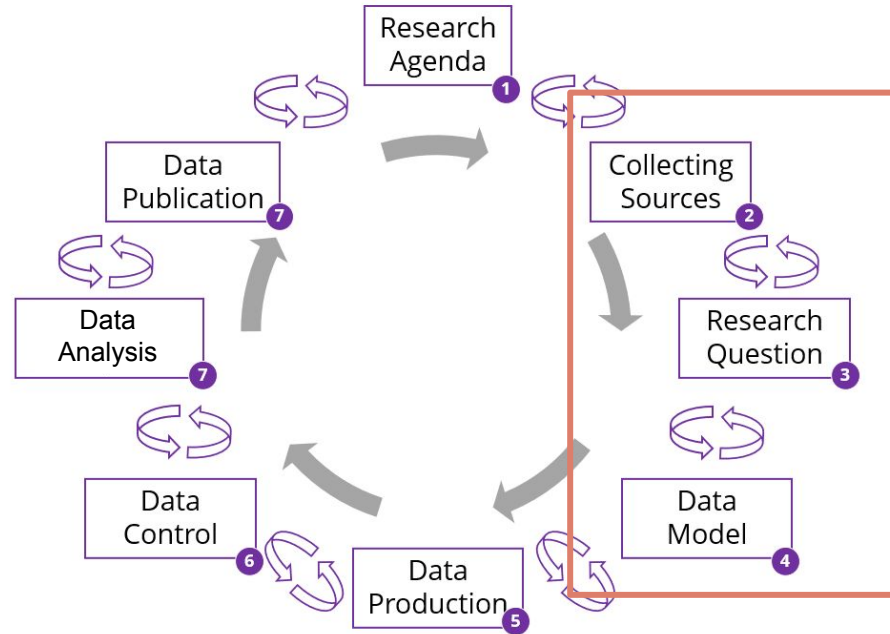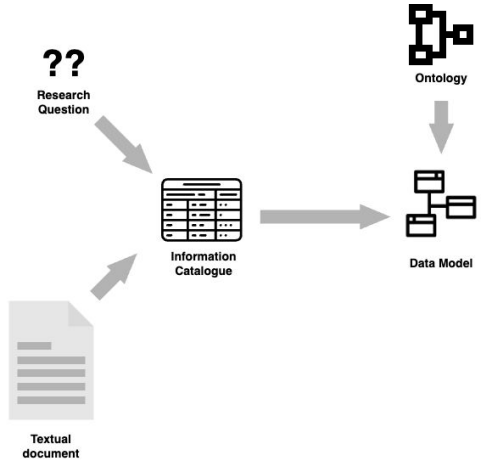  - **Our goal is to map the extracted data to CIDOC CRM and the SDHSS extension**

# Knex: The Pipeline



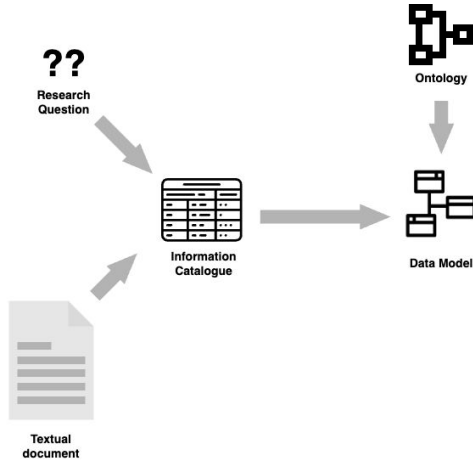https://github.com/geovistory/knex

# Developing the Data Model: Content Analysis
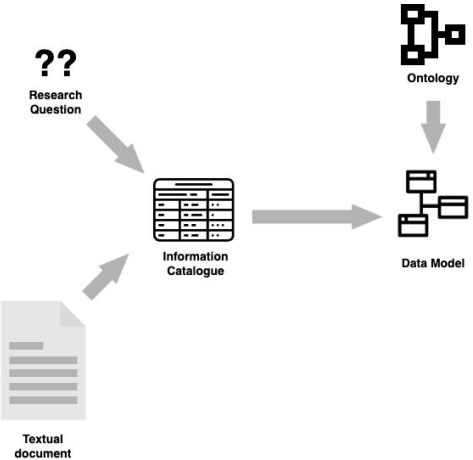
# Developing the Data Model: Content Analysis



**Theodor Tobler**

Theodor Tobler was born on 24.1.1876 in Bern and died on 4.5.1941 in Bern, prot., from Lutzenberg. Son of Johann Jacob (known as Jean), confectioner, and Adeline Lorenz, née Baumann. Married in 1903 Theda Born, daughter of Emil, architect (divorced in 1919), and married in 1919 Bertha Eschmann, daughter of Heinrich. Commercial training in Geneva and Venice. Tobler joined his father's business in 1894, adding a chocolate factory in 1899. In 1908, he created Toblerone. After leaving the company (1933), which was in need of restructuring, he bought the Klameth industrial confectionery in Bern in 1934. In 1937, he founded Typon AG in Burgdorf, which produced films for the graphics industry, and made a name for itself with new products and original advertising strategies. Member of the Bernese Masonic lodge A l'Espérance (from 1902). An entrepreneur and advocate of social reform, Tobler was also active in the pacifist and pan-European movements.
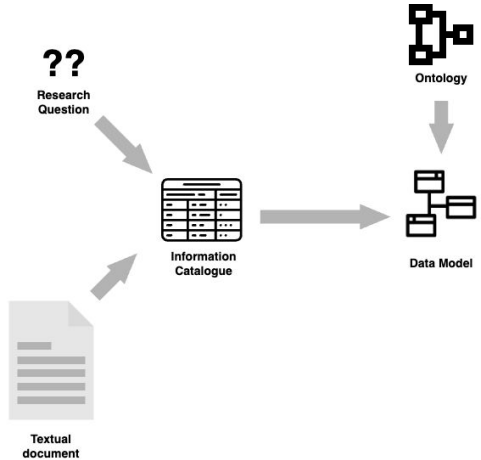
# Developing the Data Model: Information Catalogue



| Entity | Attributes |
|---|---|
| Person | Name<br>has birth place -> Place<br>Birth Date<br>has father -> Person<br>has mother -> Person<br>studied -> Study<br>has occupation -> Occupation |
| Study | has study place -> Place<br>Study Date<br>Study Domain |
| Occupation | has occupation place -> Place<br>Occupation Date<br>Related Group |
| Group | Group Name<br>has founder -> Person |
| Place | Place Name |

# Developing the Data Model: Ontology



**??** Research Question

Textual document
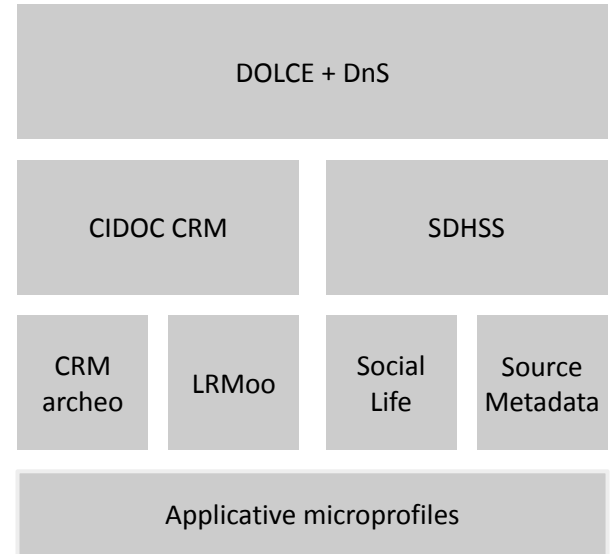
→ Information Catalogue

Ontology → Data Model

Foundational ontologies and modelling methodologies

↓

Generic, domain related core ontology

↓

Demain related extensions

↓

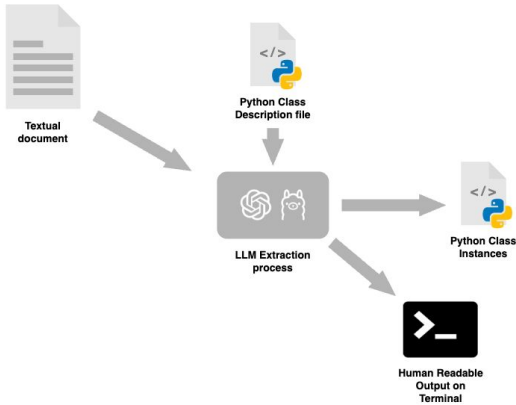Modular data models

| DOLCE + DnS | |
|---|---|
| CIDOC CRM | SDHSS |

| CRM archeo | LRMoo | Social Life | Source Metadata |
|---|---|---|---|

| Applicative microprofiles |
|---|

# Developing the Data Model: Data Model

# Knowledge Extraction Process



Textual document

Python Class Description file

LLM Extraction process

Python Class Instances

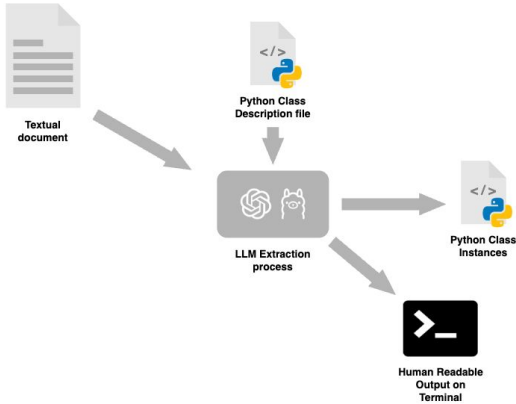Human Readable Output on Terminal

LLMs struggle to understand the event-centric model of CIDOC CRM

The solution adopted is to rely on **shortcuts**, that would then be transformed into the extended CIDOC CRM paths later on in the pipeline.

LLM steps:

- Identify each entity in the text
- For each entity, parse the text for extracting the shortcut assertions
- Store the assertions as Python class instances

# Knowledge Extraction Process: Class Description



Textual document

Python Class Description file

LLM Extraction process

Python Class Instances

Human Readable Output on Terminal

```python
# This class is to be given to a LLM.
# Classes descriptions, attributes names, fields description should be set thoroughly.
class Person(BaseModel):
    """
    Information from the text about a person.
    """

    # General informations
    name: Optional[str] = Field(default=None, description="the name of the person")
    gender: Optional[str] = Field(default=None, description='"male" or "female"')
    origins: Optional[str] = Field(default=None, description="the person origin: geographical place")
    religion: Optional[str] = Field(default=None, description="the person religion name")

    # Birth
    birth_date: Optional[str] = Field(default=None, description="the person birth date")
    birth_place: Optional[str] = Field(default=None, description="the birth place of the person (geographical place)")

    # Death
    death_date: Optional[str] = Field(default=None, description="the person death date")
    death_place: Optional[str] = Field(default=None, description="the birth place of the person (geographical place)")

    # Genealogy
    father_name: Optional[str] = Field(default=None, description="the father of the person")
    mother_name: Optional[str] = Field(default=None, description="the mother of the person")
```
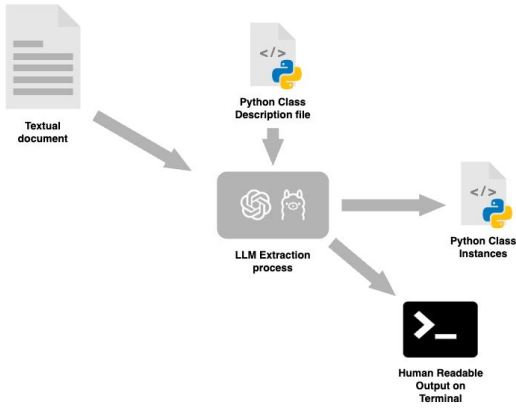
# Knowledge Extraction Process: Output



```
Persons found in the text: Theodor Tobler, Johann Jacob (Jean), Adeline Lorenz, Theda Born, Emil, Bertha Eschmann, Heinrich

==== Extracting information about: Theodor Tobler ====
name: Theodor Tobler
birth_date: 1876.01.24
birth_place: Bern
death_date: 1941.05.04
death_place: Bern
father_name: Johann Jacob (known as Jean)
mother_name: Adeline Lorenz, née Baumann

==== Extracting information about: Johann Jacob (Jean) ====
name: Johann Jacob (known as Jean)

==== Extracting information about: Adeline Lorenz ====
name: Adeline Lorenz

==== Extracting information about: Theda Born ====
name: Theda Born
father_name: Emil

==== Extracting information about: Emil ====
name: Emil

==== Extracting information about: Bertha Eschmann ====
name: Bertha Eschmann
father_name: Heinrich

==== Extracting information about: Heinrich ====
name: Heinrich

=== Extracting activities of: Theodor Tobler ===
activities 0:
    person_name: Theodor Tobler
    place: Geneva and Venice
    activity_type: formation
    discipline: Commercial training
activities 1:
    person_name: Theodor Tobler
    date_begin: 1894.00.00
    date_end: 1933.00.00
    activity_type: job
    institution: his father's business
activities 2:
```

# Data Transformation Process: Conversion Script



Data Model

Python Transformation Script

Graph File (.ttl)

Python Class Instances

```python
def person_to_graph(person: Person, graph: Graph) -> None:
    """
    Transform an object instance into a list of entities and statements.

    Args:
        person (Person): the person to be added to the graph.
        graph (Graph): the graph to add the person to.
    """

    # If the person does not have a name, set a default value
    if not person.name:
        person.name = "Unknown Person " + str(graph.get_current_index())

    # Create the person
    person_ent = graph.create_entity_aial(c.E21_person, person.name)

    # Gender
    if person.gender:
        gender = graph.create_entity_aial(          (variable) P23_hasGender: int
        graph.create_triple(person_ent, p.P23_hasGender, gender)

    # Origins
    if person.origins:
        geoplace = graph.create_entity_aial(c.C13_geographicalPlace, person.origins)
        graph.create_triple(person_ent, p.P24_hasItsOriginsIn, geoplace)

    # Religion
    if person.religion:
        religious_identity = graph.create_entity_aial(c.C23_religiousIdentity, person.religion)
        graph.create_triple(religious_identity, p.P36_pertainsTo, person_ent)

    # Birth date
```

# Data Transformation Process: Generated Graph

# The Use Case and Demo

# The sources

Historical Dictionnary of Switzerland:

https://hls-dhs-dss.ch/

Documenting 25394 people and 2559 families
in the form of textual descriptive biographies.

The case of Theodor Tobler, inventor of the
*toblerone*.



DE **FR** IT

## Theodor Tobler

Version du: 18.12.2013

Auteure/Auteur: **Christian Lüthi** | Traduction: **Florence Piguet**

﹡ 24.1.1876 à Berne, † 4.5.1941 à Berne, prot., de Lutzenberg. Fils de Johann Jacob (dit Jean), confiseur, et d'Adeline Lorenz, née Baumann. ∞ 1) 1903 Theda Born, fille d'Emil, architecte (divorce en 1919), 2) 1919 Bertha Eschmann, fille de Heinrich. Formation de commerce à Genève et Venise. T. entra dans l'entreprise paternelle en 1894, à laquelle il adjoignit une fabrique de chocolat (1899). En 1908, il créa le Toblerone. Après avoir quitté l'entreprise (1933), qui nécessitait un assainissement, il acheta en 1934 la confiserie industrielle Klameth à Berne. En 1937, il fonda à Berthoud la société Typon AG, qui produisait des films pour l'industrie graphique, et se distingua grâce à de nouveaux produits et à des stratégies publicitaires originales. Membre de la loge maçonnique bernoise A l'Espérance (dès 1902). Entrepreneur partisan de réformes sociales, T. milita aussi dans les mouvements pacifiste et paneuropéen.

# The case of Theodor Tobler

Theodor Tobler was born on 24.1.1876 in Bern and died on 4.5.1941 in Bern, prot. of Lutzenberg. Son of Johann Jacob (known as Jean), confectioner, and Adeline Lorenz, née Baumann. 1) 1903 Theda Born, daughter of Emil, architect (divorced in 1919), 2) 1919 Bertha Eschmann, daughter of Heinrich. Commercial training in Geneva and Venice. T. joined his father's business in 1894, adding a chocolate factory in 1899. In 1908, he created Toblerone. After leaving the company (1933), which was in need of restructuring, he bought the Klameth industrial confectionery in Bern in 1934. In 1937, he founded Typon AG in Burgdorf, which produced films for the graphics industry, and made a name for itself with new products and original advertising strategies. Member of the Bernese Masonic lodge A l'Espérance (from 1902). An entrepreneur and advocate of social reform, T. was also active in the pacifist and pan-European movements.
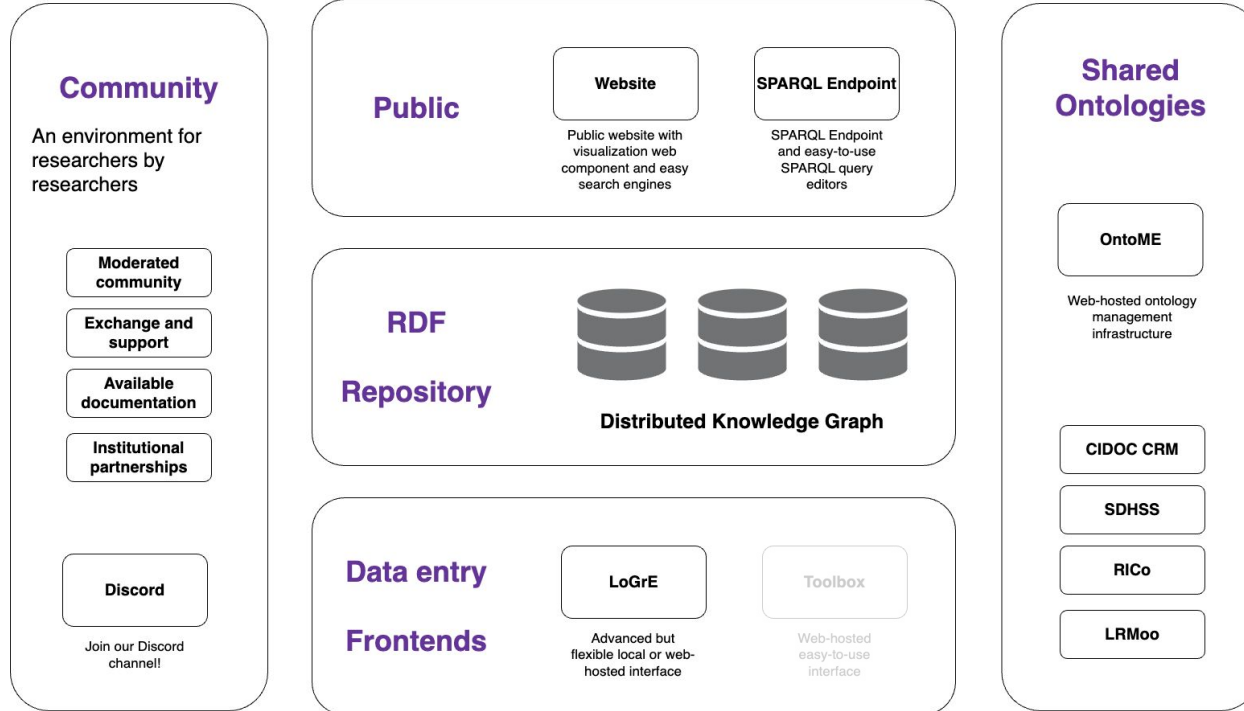
# The Geovistory Environment

https://www.geovistory.org

# The Geovistory Team

# The Geovistory Environment

Conclusions

# Challenges

- The risk of hallucination is rather low
    - It's mitigated with the class description before transforming the assertions into triples
    - Reduced by the parsing of the text for each entity identified (trying to create all the assertions at the same time created a lot of errors), but it makes the process resource intensive
- The creation of the Python Class Description file is done manually and adapted for each data transformation process:
    - This allows the creation of a more precise Knowledge Graph, especially in line with the research questions, but requires a lot of resources on our side.
- Close monitoring on the evolution of the technologies is needed as they evolve very quickly

# Next Steps

1. Make Knex more generic and reusable:
   - In order to make Knex more reusable beyond our team, we need develop process to generate Class Description files, such as with the use of SHACL ontological profiles (managed in the OntoME platform)
2. Find more efficient LLM processes:
   - Most of the LLM resources go in the multiple parsings of the text for generating the class instances. Finding ways to parse the text only once would significantly reduce the costs.
3. Integrate Knex in the graph management tool Logre:
   - Develop a more flexible and reusable GUI